

# Justification of Automated Decision-Making: Medical Explanations as Medical Arguments

Ravi D. Shankar, MS, and Mark A. Musen, MD, PhD  
Stanford Medical Informatics, Stanford University School of Medicine,  
Stanford, CA 94305-5479

*People use arguments to justify their claims. Computer systems use explanations to justify their conclusions. We are developing WOZ, an explanation framework that justifies the conclusions of a clinical decision-support system. WOZ's central component is the explanation strategy that decides what information justifies a claim. The strategy uses Toulmin's argument structure to define pieces of information and to orchestrate their presentation. WOZ uses explicit models that abstract the core aspects of the framework such as the explanation strategy. In this paper, we present the use of arguments, the modeling of explanations, and the explanation process used in WOZ. WOZ exploits the wealth of naturally occurring arguments, and thus can generate convincing medical explanations.*

## 1 THE PROPOSAL – EXPLANATIONS AS ARGUMENTS

In conversation, humans use arguments to support beliefs and claims by providing evidences. Lawyers argue their clients' cases, writers propound their beliefs, and physicians justify their diagnoses. Explanation is the way that a knowledge-based system similarly justifies its conclusions to its users. As requested by the user, the system presents different levels of information in support of its claims.

We are developing WOZ, a multi-agent framework [1] that provides explanations by using arguments. WOZ explains the claims of EON [2], a knowledge-based system architecture that provides physicians with decision-support in protocol-based care. EON's problem-solving components use explicit models of medical-domain and clinical-protocol knowledge. WOZ, like EON, uses explicit models that abstract the explanation strategy and the agent architecture. The explanation strategy defines what information WOZ uses to justify an EON claim. The strategy model uses an argument structure proposed by Toulmin in his theory of reasoning [3]. Toulmin's structure allows the recipient of the argument to identify the different elements needed to support the claim.

In this paper, we explain how WOZ does explanation modeling based on arguments. We describe EON and WOZ (Section 2), then discuss what the general role of arguments is and how arguments relate to explanations (Section 3.1). We introduce Toulmin's argument structure and summarize its use in previous explanation frameworks (Section 3.2). We then describe our

explanation models and their use in the explanation process (Section 3.3). Finally, we discuss the merits of modeling explanations as arguments, and show how the models might be extended to make explanations more persuasive (Section 4).

## 2 THE GROUNDS – EON AND WOZ

A clinical protocol such as the HIV/AIDS protocol CBCT-P001 has a set of eligibility criteria that clinicians use to decide whether a patient can be treated in accordance with that protocol. One of the decision-support systems that we have built for use within the EON architecture determines a patient's eligibility for a given protocol [4]. The system contains a set of general components that work together to automate the eligibility-determination task (Figure 1, shaded area). The database-mediator component performs temporal database management and data abstractions on patient data. The domain-model component provides the domain-specific terms and relations used by the system. The protocol

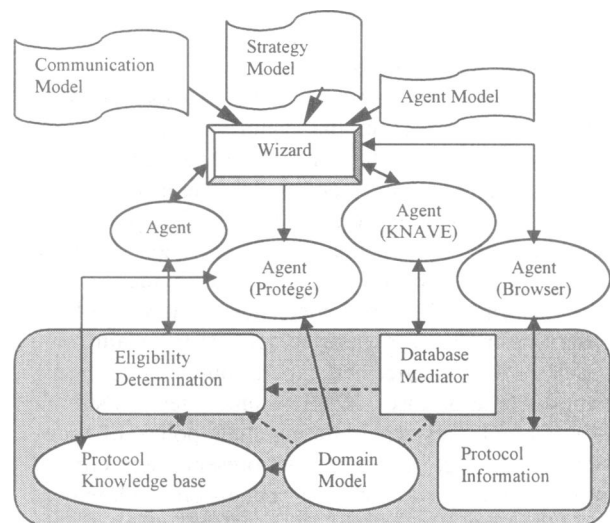


Figure 1: EON and WOZ. The EON components are shown inside the gray box. These components interact with each other to provide decision support. WOZ consists of a set of visualization agents, one agent per component. Each agent has a set of graphical user interfaces. The Wizard controls the agent-agent interactions. It takes in as input the communication model, the agent model, and the strategy model. The Wizard parses an explanatory query, consults the explanation strategy to decide what information should be presented as an explanation, and directs the appropriate agents to present that information.

knowledge-base component provides structured protocol knowledge used by the eligibility-determination component. The protocol-information component displays for the end user texts of protocol documents.

WOZ employs a collection of cooperating visualization agents (Figure 1) that are closely associated with the components of EON. The agents handle the higher-level aspects of explanation, such as presentation of information and interaction with the end user. The components provide the lower-level data that drive the explanation, such as the raw patient data and the reasoning knowledge. An explanation engine, the Wizard, uses the explanation strategy model to decide which agents should present what information in response to a user query. The visualization agents, the Wizard, the explanation strategy, the agent interactions, and the components constitute the explanation of the whole system.

### 3 EXPLANATION MODELS USING ARGUMENT STRUCTURES

The core aspect of our explanation framework is the explanation strategy that defines what constitutes an explanation for a claim. It identifies the different components of an explanation such as the claim itself, the medical evidences that support the claim, the strength of the claim, the evidences that contradict the claim, and the patient data that EON used in determining the claim. This explanation strategy is like that used by people who use arguments to support a claim.

#### 3.1 About Arguments

A person states a claim to others, and then uses arguments to increase their belief in the claim. These arguments include presenting evidences related to the claim, generally using a structured format. Arguments are widely used by lawyers and writers. A lawyer first presents a claim in her opening statement at a trial. Then, to prove the claim, she describes the physical evidence that supports the claim. She then identifies the laws apply to the case. She may call experts to support her claim. She may back up the evidence and expert opinions by referring to historical data or cases. Writers sometimes use similar structured arguments when stating their points of view. Physicians may use arguments to present their clinical diagnoses, and may support their conclusions with patient data, medical knowledge, and other related cases. Given this role of arguments in human-human interactions, we can see how arguments can be synonymous with explanations in human-system interactions.

#### 3.2 Toulmin's Argument Structure and Explanations

We used Toulmin's argument structure to identify, organize and present information related to the explanation strategy. Toulmin, a philosopher, has

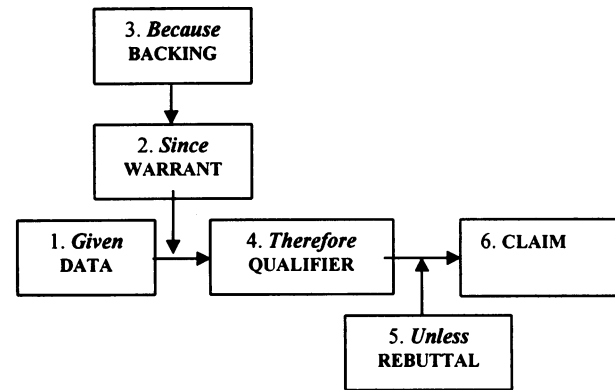


Figure 2: Toulmin's Argument Structure. The structure reads given Data, therefore Claim, since Warrant, because Backing, unless Rebuttal. The elements of the structure and the relationship among them can be used to generate explanations.

developed a pragmatic method of reasoning [3]. Central to his logic is his six-element argument structure (Figure 2) by which claims can be argued, regardless of the context of the argument:

1. **Data:** The particular facts about a situation on which a claim is made
2. **Warrant:** The knowledge that justifies a claim made using the data
3. **Backing:** The general body of information or experience that validate the warrant
4. **Qualifier:** The phrase that shows the confidence with which the claim is supported to be true
5. **Rebuttal:** The anomaly that shows the claim not to be true
6. **Claim:** The assertion or conclusion put forward for general acceptance

Wick [5] pointed out how early research in explanation has, without stated intent, evolved to engulf Toulmin's argument structure. Ye [6] used Toulmin's argument structure to study the value of explanation in expert systems for auditing. Ramberg [7] describes a multiple-explanation construction model that constructs explanations for an expert system in the domain of protein purification.

#### 3.3 Explanation Models of WOZ

The explanation space of EON comprises the patient data (e.g., laboratory parameters), the medical domain knowledge (e.g., AIDS domain concepts), the clinical-protocol knowledge (e.g., eligibility criteria for protocol CBCT-P001), the reasoning knowledge (information on the reasoning process), and other relevant external information (e.g., clinical-protocol text). The explanation process in the WOZ framework includes (1) identifying the distinct elements of the explanation space that are required to satisfy user's explanatory query, (2) obtaining the required information from the appropriate agents, and

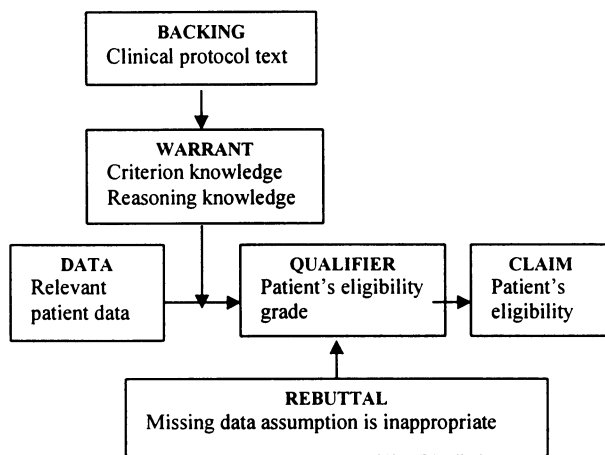


Figure 3. A Meta-argument for a class of claims. The structure illustrated here is for the class *patient's eligibility*. The meta-argument is specified explicitly and stored in the explanation strategy knowledgebase. Note that, in this example, all the elements of the Toulmin structure have been specified. Typically, the evidences available to support a claim drives what elements get filled.

(3) presenting the explanation in a coherent manner. To aid in this explanation process, we have designed three declarative models, the Strategy Model, the Communication Model and the Agent Model. These explanation models abstract the core elements of the explanation process: the explanation strategy, the agent services, and the terminology used in the strategy and by the agents. We built the explanation models using Protégé, a software system that who uses to create knowledge models. We also generated the corresponding knowledge-acquisition tools using Protégé.

### 3.3.1 The Strategy Model

We modeled the explanation strategy as *arguments*, using Toulmin's argument structure. The elements of the argument structure for a claim identify the information needed for explanation of that claim. Appropriate WOZ explanation agents provide the needed information, and the explanation strategy can use the relationships among the elements of the argument structure to present the explanation consistently and clearly. We define two types of arguments in our explanation strategy, (1) the meta-argument and (2) the concrete argument.

A **meta-argument** conceptualizes arguments for a *class of claims*. We state explicitly the elements of the meta-argument structure with abstract descriptions of appropriate pieces of information in EON's explanation space. We illustrate a meta-argument for explaining patient eligibility score in Figure 3. The meta-argument is "The criterion knowledge and the reasoning knowledge support the graded eligibility score determined using the data; the clinical-protocol text provides the basis for the criterion knowledge and the reasoning knowledge; the

claim is not true if the missing data assumption, if any, is inappropriate."

A **concrete argument** defines an argument for a *specific claim* in a class of claims. It follows that a concrete argument is an instance of a meta-argument. We specified a meta-argument for explaining patient eligibility score (Figure 3). From this meta-argument, we can derive a concrete argument for explaining patient's eligibility score for a particular criterion such as the *platelet-count criterion* (Figure 4). The abstract descriptions of the pieces of information in the meta-argument are substituted with the actual information related to the computation of this criterion's eligibility score. With this concrete argument, WOZ can generate an explanation to support the criterion's eligibility score.

Meta-arguments and concrete arguments have a class-instance relationship. The explanation-strategy knowledge consists of meta-arguments, one for each class of claims the system makes. Since the explanation strategy is modeled explicitly, we can acquire the meta-arguments from the experts using a knowledge-acquisition tool. The concrete arguments are derived from these meta-arguments at runtime during the explanation process.

### 3.3.2 The Communication Model

There are concepts and terms that are used when developers construct the meta-arguments and when WOZ agents communicate with one another (see Figure 1). The communication model defines the following concepts:

- Action verbs that are mainly used in agent-agent

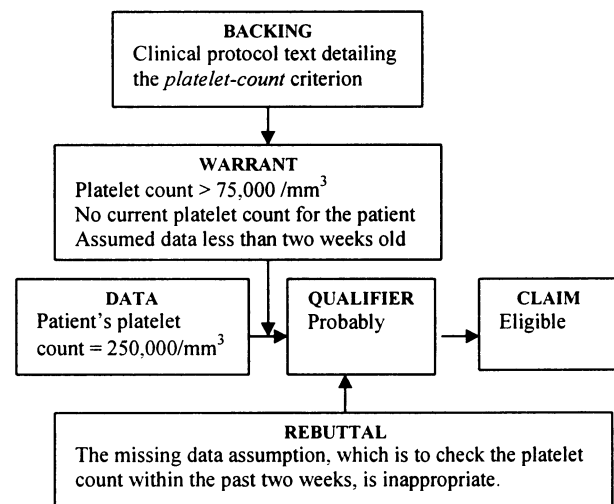


Figure 4. A Concrete argument for a specific claim. The structure illustrated here is for the claim *patient's eligibility for platelet-count criterion*. A justification for the claim can be made using the elements of this concrete argument structure.

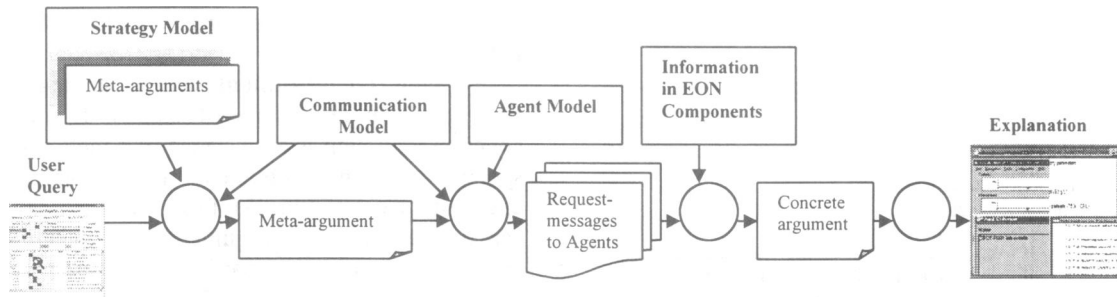


Figure 5. The WOZ Explanation Process. The user submits query and receives explanation using GUIs. The numbered circles refer to the following action points: 1) Wizard selects appropriate meta-argument; 2) Wizard asks appropriate agents to provide information specified in the meta-argument; 3) WOZ Agents provide information; 4) Wizard generates explanation..

interactions (*explain* and *show*).

- Concepts related to the problem domain (*protocol identifier*, *criterion identifier*, and *eligibility score* in the patient eligibility determination domain).
- Concepts and terms related to the particular area of medicine, such as breast cancer or HIV/AIDS, to which the knowledge-based system is applied (*clinical findings* and *medical interventions*).

### 3.3.3 The Agent Model

The visualization agent in the WOZ architecture provides information in the EON component that it is representing (Figure 1). The agent presents this information using visual media such as text, graphics and video. The agent model encapsulates the characteristics of the agents, such as the agent's identifier, and the services that the agent provides. The definition of a *service* includes the required inputs, the type of information provided, and the presentation medium. An example instantiation of the agent model is information on an agent identified as *KNAVE* (Figure 1) that presents *patient data* graphically, taking *patient identifier* as input.

### 3.4 A Dialog

To explain a specific claim, at runtime, WOZ selects the appropriate meta-argument in the strategy knowledge base, identifies the explanation information, obtains the information from appropriate agents, derives the corresponding concrete argument, and generates the explanation by organizing the presentation of the agents (Figure 5). We demonstrate this explanation in a visual dialog between WOZ and a user of EON. The user submits queries and receives explanations using direct manipulation via graphical user interfaces. This dialog is based on the examples that we used previously to illustrate the meta-argument and the concrete argument.

1. EON has determined the patient's eligibility for protocol CBCT-P001, and displays the eligibility scores of the protocol overall and of the protocol's individual criteria. Now the user is interested in the details of a specific criterion of the protocol.

2. User submits a request to display the patient's data used to compute the eligibility score for the *platelet-count* criterion.
3. **Wizard Agent** (Figure 1) receives the details of the request. It recognizes that the eligibility score of the *platelet-count* criterion belongs to the class of eligibility scores. It selects the meta-argument of the eligibility scores (Figure 3) from the strategy knowledge base. It then consults the agent knowledge base and requests the appropriate agents to provide the information pieces identified in the meta-argument.
4. **Agents** fill in the realtime values for the *platelet-count* criterion to get the concrete argument for the criterion (Figure 4). For example, the *KNAVE* agent fills the *data* element, the *Protégé* agent fills the *warrant* element, and the *Browser* agent fills the *backing* element. The agents also generate visual presentations of the elements.
5. **WOZ** uses the *data* element to answer the user. It displays the *KNAVE* agent's presentation that the patient had a platelet count = 250,000/mm<sup>3</sup>.
6. User requests to explain the eligibility score for the platelet-count criterion.
7. **WOZ** taps into the *warrant* element to answer the user. It displays the *Protégé* agent's presentations that the eligibility criterion is "platelet count should be > 75,000/mm<sup>3</sup>", and the missing data assumption is "assume that the platelet count is valid if less than 2 weeks old, and that the score is *probably eligible* if within the range."
8. User requests more information on the eligibility criterion.
9. **WOZ** directs the user to the information in the *backing* element. It displays the *Browser* Agent's presentation of the protocol text that describes the *platelet-count* criterion. This presentation may also indicate the intentions of the protocol authors behind the *platelet-count* criterion.

This dialog demonstrates how the concrete argument structure aids WOZ in deciding what information to present when.

## 4 DISCUSSION

We have demonstrated how WOZ can justify claims made by knowledge-based systems such as EON by using arguments. We used Toulmin's argument structure to integrate and present explanation information from varied sources: the EON components. Acquiring explanation strategies for particular domains includes identifying the classes of claims that the system need to make and defining one meta-argument for each class. We believe that the number of classes is within reasonable limits for a component-based system such as EON. For example, the EON patient-eligibility-determination system requires only one class of claim (the patient's eligibility score).

The dialog that we presented in Section 3.4 could have taken place between a physician and a patient. Horton [8] and Dickinson [9] proposed that physicians use Toulmin's argument structure to organize medical evidences supporting their diagnoses or treatment plans. These proposals support our contention that it is appropriate to use an argument approach to provide medical explanations.

There are explicit relationships among the elements of Toulmin's argument structure. We can use these relationships when presenting explanations. Other relationships in the argument structure, however, may have to be defined and considered in the presentation. For example, when there is more than one item in any element of the argument – say the *Warrant* – in what order do we present these items? How do we expose the relationships among these items? One method is to employ the Rhetorical Structure Theory (RST) [10] that was developed mainly for text analysis and text generation. RST maintains that, in most coherent discourse, consecutive discourse elements are related by a small set of rhetorical relations that is defined by the theory. Three of these relations are *Condition*, *Elaboration* and *Sequence*. Many natural-language-generation systems rely on the rhetorical relations defined in RST. We can use RST by describing various rhetorical relations among the different items of the warrant. By providing such connectives, we can strengthen the cohesiveness of the multiple but related items, thereby enhancing the clarity in the presentation.

There is always a question of how to tailor explanations to the user in a way that will enhance the user's acceptance of the claims. One approach is to extend the argument structure by defining multiple subarguments. We can envision this structure as multiple argument structures having the same *Claim*, *Data*, and *Modifier* but possibly different *Warrant*, *Backing* and *Rebuttal*. A superimposition of these structures will result in an argument structure that contains multiple subarguments. For the same *claim* and *data*, we can then create many subarguments each providing a different flavor of the same argument. When WOZ is providing an explanation to a user, WOZ can employ a suitable

subargument, thus providing tailored explanations. Naturally, this design presupposes the existence of a user model that abstracts the user profile and preferences.

## 5 CONCLUSION

We showed how medical explanations could be expressed as medical arguments. Our explanation strategy uses a widely recognized argument structure, and can mirror naturally occurring medical arguments. Thus, our approach can generate convincing medical explanations.

### Acknowledgements

This work has been supported, in part, by grant LM05708 from the National Library of Medicine.

We thank Lyn Dupre for her valuable editorial comments.

### References

1. Shankar RD, Tu SW, Musen MA. A declarative explanation framework that uses a collection of visualization agents, *Journal of the American Medical Informatics Association*, 1998; symposium supplement, 602–606.
2. Musen MA, Tu SW, Das AK, Shahar Y. EON: A component-based approach to automation of protocol-directed therapy, *Journal of the American Medical Informatics Association*, 1996; 3(6), 367–388.
3. Toulmin S. The uses of argument, *Cambridge University Press*, Cambridge MA, 1958.
4. Tu SW, Kemper CA, Lane NM, Carlson RW, Musen MA. A methodology for determining patients' eligibility for clinical trials, *Methods of Information in Medicine*, 1993; 32(4), 317–325.
5. Wick MR. Expert system explanation in retrospect: a case study in the evolution of expert system explanation, *Journal of Systems and Software*, 1992; 19(2), 159–169.
6. Ye LR. The value of explanation in expert systems for auditing: An experimental investigation, *Expert Systems with Applications*, 1995; 6(4), 543–556.
7. Ramberg R. Construing and testing explanations in a complex domain, *Computers in Human Behavior*, 1996; 12(1), 29–48.
8. Horton R. The grammar of interpretive medicine, *Canadian Medical Association Journal*, 1998; 158, 245–249.
9. Dickinson, HD. Evidence-based decision-making: an argumentative approach, *International Journal of Medical Informatics*, 1998; 51, 71–81.
10. Mann WC, Thompson S. Rhetorical structure theory: A theory of text organizations, *Information Sciences Institute Technical Report Number RS-87-190*, University of Southern California, Marina del Rey, CA, 1987.